

---

# CREDIBLE, UNRELIABLE OR LEAKED?: EVIDENCE VERIFICATION FOR ENHANCED AUTOMATED FACT-CHECKING

---

Zacharias Chrysidis<sup>1</sup>, Stefanos-Iordanis Papadopoulos<sup>1,2</sup>, Symeon Papadopoulos<sup>2</sup>, and Panagiotis C. Petrantonakis<sup>1</sup>

<sup>1</sup>Department of Electrical & Computer Engineering, Aristotle University of Thessaloniki.

<sup>2</sup>Information Technology Institute, Centre for Research & Technology, Hellas.  
zachoschrysidis@gmail.com, {stefpapad,papadop}@iti.gr, ppetrant@ece.auth.gr

## ABSTRACT

Automated fact-checking (AFC) is garnering increasing attention by researchers aiming to help fact-checkers combat the increasing spread of misinformation online. While many existing AFC methods incorporate external information from the Web to help examine the veracity of claims, they often overlook the importance of verifying the source and quality of collected “evidence”. One overlooked challenge involves the reliance on “leaked evidence”, information gathered directly from fact-checking websites and used to train AFC systems, resulting in an unrealistic setting for early misinformation detection. Similarly, the inclusion of information from unreliable sources can undermine the effectiveness of AFC systems. To address these challenges, we present a comprehensive approach to evidence verification and filtering. We create the “CREDible, Unreliable or LEaked” (CREDULE) dataset, which consists of 91,632 articles classified as Credible, Unreliable and Fact-checked (Leaked). Additionally, we introduce the Evidence VERification Network (EVVER-Net), trained on CREDULE to detect leaked and unreliable evidence in both short and long texts. EVVER-Net can be used to filter evidence collected from the Web, thus enhancing the robustness of end-to-end AFC systems. We experiment with various language models and show that EVVER-Net can demonstrate impressive performance of up to 91.5% and 94.4% accuracy, while leveraging domain credibility scores along with short or long texts, respectively. Finally, we assess the evidence provided by widely-used fact-checking datasets including LIAR-PLUS, MOCHEG, FACTIFY, NewsCLIPPings+ and VERITE, some of which exhibit concerning rates of leaked and unreliable evidence.

**Keywords** Deep Learning, Misinformation Detection, Automated Fact-Checking, Evidence Filtering, Information Leakage

## 1 Introduction

Misinformation has become an increasingly prevalent issue today, causing negative impacts to individuals and society [1]. With the rapid spread of online platforms and social media, fake or misleading information has become alarmingly widespread, posing significant challenges to informed decision-making and societal trust [2]. In the battle against misinformation, many fact-checking platforms such as Snopes<sup>1</sup>, PolitiFact<sup>2</sup> and Reuters<sup>3</sup> have emerged, where journalists manually review a plethora of claims sourced from news articles and social media. Nonetheless, manual fact-checking is time-consuming and can not always keep pace with the rate at which misinformation spreads. Recently, researchers in natural language processing [3], computer vision [4] and multimodal learning [5] have begun exploring Automated Fact-Checking (AFC). AFC involves tools and systems that help professional fact-checkers to combat misinformation more efficiently by automating pivotal aspects of fact-checking including claim detection, evidence retrieval and claim verification [6].

---

<sup>1</sup>Snopes: <https://www.snopes.com/>

<sup>2</sup>Politifact: <https://www.politifact.com/>

<sup>3</sup>Reuters: <https://www.reuters.com/fact-check/>

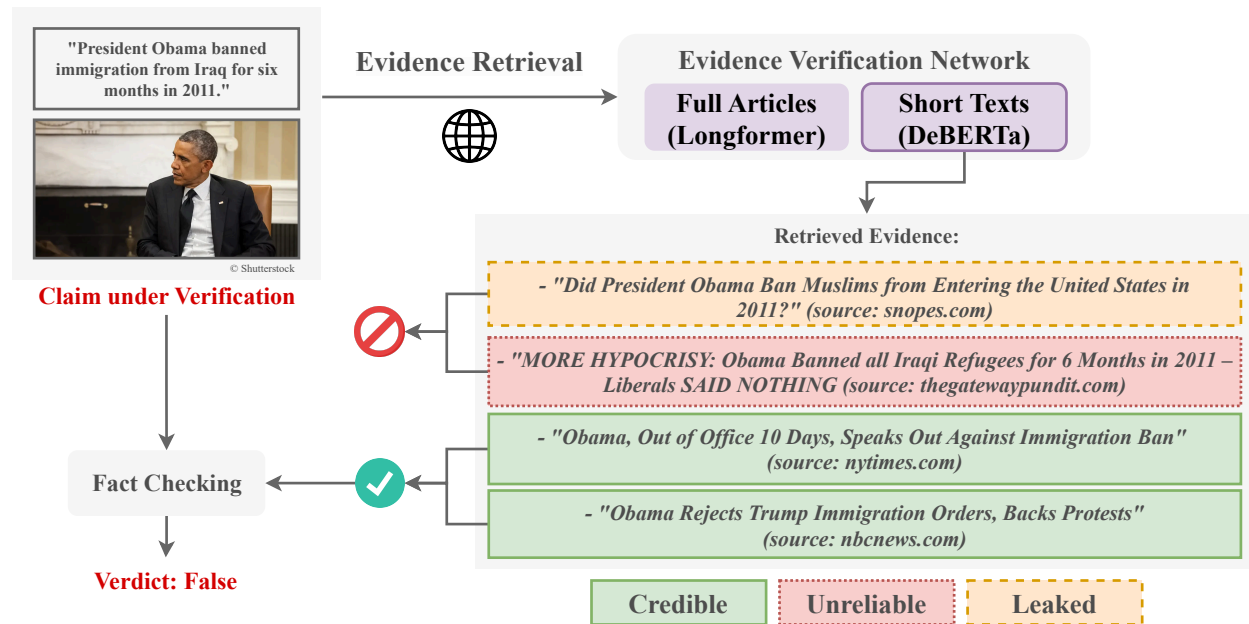


Figure 1: Pipeline of automated fact-checking leveraging the proposed Evidence Verification Network.

To further improve AFC, some systems leverage external evidence extracted from the Web using search engines. By tapping the potential of the entire Web as a knowledge source to help support or refute a claim, these models enhance their ability to verify news pieces more accurately. However, a prevalent issue in external knowledge retrieval from the Web is the lack of evidence filtering mechanisms. Inadequate or rudimentary filtering results in the inclusion of irrelevant or unreliable information, potentially compromising the accuracy of fact-checking systems. Furthermore, Glockner et al. [7] define “two requirements that the evidence in datasets must fulfill for realistic fact-checking: It must be (1) sufficient to refute the claim and (2) not leaked from existing fact-checking articles”. Otherwise the AFC model would learn to rely on previously fact-checked information when trying to detect new emerging misinformation, where fact-checks are not yet available. The problem of “leaked evidence” is quite under-researched yet crucial for realistic and effective fact-checking.

Motivated by these observations, we propose a new evidence verification and filtering approach to address the issue of leaked and unreliable evidence in AFC. Firstly, we construct the “CREDible, Unreliable or LEaked” (CREDULE) dataset, by modifying, merging, and extending MultiFC [8], Politifact [9], PUBHEALTH [10], NELA-GT [11, 12, 13, 14, 15, 16], Fake News Corpus [17], and Getting Real About Fake News [18]. These established datasets contain short texts (titles) as well as the long texts (full articles) of the news articles. We extract the article bodies, where they are not given, and other meta-data to better balance the classes. The final CREDULE dataset consists of 91,632 pieces, equally distributed in three classes: “Credible”, “Unreliable” and “Fact-checked” (or Leaked).

The goal is to develop a model capable of detecting the information that a model retrieves from the Web so as to avoid leakage and unreliable sources, as seen in Figure 1. To this end, we also propose EVVER-Net a neural network that detects leaked (fact-checked) and unreliable evidence pieces during the evidence retrieval process and only allows credible information to pass to the AFC model. We experiment with various pre-trained Transformer-based encoders for both short texts, such as DeBERTa [19], CLIP [20], T5 [21], and long texts, Long T5 [22] and Longformer [23], as well as baseline methods like Count Vectorizer and TF-IDF. Additionally, we integrate domain credibility scores from the Media Bias/Fact Check (MBFC) website <sup>4</sup> to enhance classification accuracy.

To show the efficacy and usefulness of the classifier, we examine the collected evidence in widely used AFC datasets such as the LIAR-PLUS [24], FACTIFY [25], MOCHEG [26], VERITE [27] datasets and the evidence collected by Abdelnabi et al. [28] for the NewsCLIPPings dataset; referred to as NewsCLIPPings+ for simplicity. Our analysis shows that the collected evidence often contain information leaked from fact-checking articles or provide unreliable information.

The contributions of our work can be summarized as follows:

<sup>4</sup><https://mediabiasfactcheck.com/>

- We propose a novel approach to detecting and filtering out leaked evidence and unreliable information in AFC systems.
- We construct CREDULE, a large-scale, balanced and diverse dataset comprising “Credible”, “Unreliable” and “Fact-checked” news articles<sup>5</sup>.
- We introduce EVVER-Net which demonstrates impressive performance of up to 91.5% and 94.4% accuracy on CREDULE while leveraging domain credibility scores along with short or long texts, respectively.
- We use EVVER-Net to examine the evidence of widely used fact-checking datasets, where we identify concerning rates of leaked and unreliable evidence.

## 2 Related Work

### 2.1 Fact-checking Datasets

A plethora of datasets have been curated to facilitate fact-checking tasks and train robust misinformation detection models. One widely used text-based dataset is FEVER [29]. It consists of 185,445 pieces generated by human annotators extracting claims from Wikipedia and mutating them in various ways, some of which alter their meaning. Many datasets also contain claims extracted from fact-checking websites. For instance, Alhindi et al. [24] created the LIAR-PLUS dataset, comprising 12,836 statements taken from Politifact and labeled by humans for truthfulness. The authors also automatically extracted justifications provided in the associated fact-checking articles. Others include FakeNewsNet [30], MultiFC [8], Politifact Fact Check [9] or WatClaimCheck [31] and some specialize in various domains such as politics (ClaimBuster [32] or Truth of Varying Shades [33]) or health (PUBHEALTH [10]). Additionally, researchers have created datasets containing both credible and fake news pieces. For example, the NELA-GT Datasets (2017-2022) [11, 12, 13, 14, 15, 16] are large corpora containing articles from both credible and non-credible news outlets. FakeNewsCorpus [17] contains both credible and fake news scraped from a curated list of 1001 non-credible domains. Conversely, the Getting Real about Fake News Dataset [18] extracted articles from 244 websites tagged as “bullshit” by the BS Detector Chrome Extension.

Multimodal datasets play a vital role in AFC by incorporating diverse types of information. They offer a more comprehensive representation of real-world claims, enabling fact-checking models to consider a broader range of evidence sources. MOCHEG [26] is a dataset consisting of 21,184 textual claims from Politifact and Snopes that also provides image and textual evidence collected from fact-checking articles. Focusing on social media content, Boididou et al. [34] built a dataset for the MediaEval 2016 Verifying Multimedia Use (VMU) challenge that comprises tweets and images. Another notable multimodal dataset is FACTIFY [25], containing 50,000 claims accompanied by 100,000 images. Collected from reliable US and Indian sources, as well as reputable fact-check websites, FACTIFY provides a diverse range of real-world data for fact-checking purposes. Researchers have also been experimenting with synthetically created multimodal misinformation. Aneja et al. [35] curated COSMOS, a dataset comprising 200K images with 450K textual captions from various news websites (credible and fact-check), blogs, and social media posts and randomly sampled negative “de-contextualized” samples. Similarly, the NewsCLIPPings dataset [36] comprises both pristine and convincing falsified (‘out-of-context’) image-caption pairs, providing examples of how misinformation can be spread through visual content. But instead of relying on “naive” random negative samples, the authors leverage CLIP [20] as well as Person and Scene Matching models to create “hard” negative samples. Built on the VisualNews corpus [37], NewsCLIPPings contains examples that misrepresent the context, place, or people in the image. Finally, the VERITE dataset was recently developed as an evaluation benchmark for multimodal misinformation detection and accounts for unimodal biases [27].

### 2.2 Evidence Collection

In the pursuit of enhancing the effectiveness of AFC, researchers have recognized the value of gathering and utilizing external information from the Web. Models often leverage popular search engines to access a vast repository of information that can supplement existing datasets. Popat et al. [38] utilized claims from Snopes and information about hoaxes and fictitious persons from Wikipedia to conduct their experiments, employing these claims and hoaxes as queries to the Google search engine. In subsequent work, the authors attempted to rank results based on the credibility of their sources [39]. Samarinas et al. [40] extended the FEVER dataset [29] to create the Factual-NLI+ dataset, incorporating synthetic examples and noise passages from web search results. Retrieving the top 30 results from the Bing Search engine for each claim in the FEVER dataset, they retained results with the highest BM25 score. Similarly, Abdelnabi et al. [28] collected external information from the Web to verify image-caption claims in the NewsCLIPPings

<sup>5</sup>We release our code at: <https://github.com/mever-team/credule-dataset>

[36] dataset. Employing an inverse search mode via the Google Vision APIs, they retrieved textual evidence such as text snippets and image captions and then utilized the caption as textual queries to search for images using the Google Custom Search API. More recently, a similar approach was adopted to augment the VERITE evaluation benchmark [27] with external information from the Web [41].

## 2.3 Addressing Leaked Evidence

While leveraging external evidence from the Web holds promise for enhancing the accuracy of AFC systems, it also presents challenges, particularly concerning the presence of leaked evidence. The criteria outlined by [7] stress the importance of ensuring that evidence used is not leaked from existing fact-checking articles. However, relying on search engines to retrieve external evidence often results in leaked information being included unintentionally. Even after excluding search results that point to the claim’s fact-checking article, leaked evidence persists. This can occur when different organizations verify the same claims or disseminate fact-checkers’ verifications. Khan et al. [31] also highlight this issue, noting that ‘premise’ articles may indirectly leak the veracity label. Glockner et al. [7] express doubts on whether automated approaches can realistically refute harmful real-world misinformation as many of existing approaches fail to overcome the information leakage problem. This underscores the need for robust mechanisms to filter out leaked evidence and maintain the credibility of AFC processes.

## 2.4 Evidence Filtering

Filtering and verifying external evidence pose significant challenges for AFC systems. Many existing models lack robust filtering mechanisms or rely on rudimentary approaches, potentially resulting in the inclusion of irrelevant, untrustworthy, or leaked evidence. Addressing this challenge requires the development of advanced filtering techniques capable of discerning reliable sources and accurate information. For example, Abdelnabi et al. [28] implemented a filtering method that discards evidence items matching the query and originating from the same website, utilizing techniques such as removing punctuation and converting captions to lowercase for textual evidence and employing perceptual hashing for images. Karadzov et al. [42], focusing on the trustworthiness of the source, disregarded evidence from domains considered unreliable based on manual checks of the most frequent domains in search results. However, this approach may not effectively identify unreliable sources that appear less frequently. Popat et al. [38] adopted an approach to assess the reliability of Web sources by determining the AlexaRank and PageRank of each source. AlexaRank measures website popularity based on its unique visitors and page views, while PageRank assesses website importance by considering the number and quality of links to and from the website. While such approaches provide valuable insights into the authority and popularity of Web sources, they may not accurately reflect their credibility from a fact-checking standpoint. In their recent study, Schlichtkrull et al. [43], employed a custom Google search tool to mitigate temporal leaks, ensuring that only documents published before the claim date were retrieved. However, this approach can introduce noise, as the dates provided by Google Search are not consistently accurate.

# 3 METHODOLOGY

## 3.1 Problem Formulation

The AFC process typically unfolds through three sequential stages: claim detection, evidence retrieval, and claim verification, as outlined by Guo et al. [6]. However, a fundamental challenge arises in the evidence retrieval stage, as not all available information can be considered trustworthy. Guo et al. [6] underscore this issue, highlighting that reliance on single authoritative sources may overlook contradicting or untrustworthy evidence. Additionally, Glockner et al. [7] emphasize the problem of information leakage, where existing AFC models incorporate leaked evidence, compromising the integrity and realism of the fact-checking process.

In response to these challenges, we construct CREDULE, a large scale dataset containing news articles from various sources. These articles are classified in three classes: “Credible”, “Unreliable” and “Fact-checked”. We also create the Evidence VERification Network (EVVER-Net) and train it on CREDULE, with the aim to detect leaked and unreliable evidence and only allow credible information to be examined by an AFC system. Specifically, the EVVER-Net  $\mathcal{G}(\cdot)$  aims to classify each textual evidence snippet  $e_i$  as fact-checked(0), credible(1), or unreliable(2), denoted by  $\mathcal{G}(e_i)$ . Let  $C$  denote the claim under scrutiny,  $E_T = \{e_1, e_2, \dots, e_N\}$  represents the set of  $N$  textual evidence pieces and  $E_I$  symbolizes the collection of image evidence related to  $C$ . After applying the classifier  $\mathcal{G}(\cdot)$ , we obtain the filtered set of textual evidence pieces  $E'_T$ , denoted as:

$$E'_T = \{e_i \mid \forall e_i \in E_T, \mathcal{G}(e_i) = 1\} \quad (1)$$

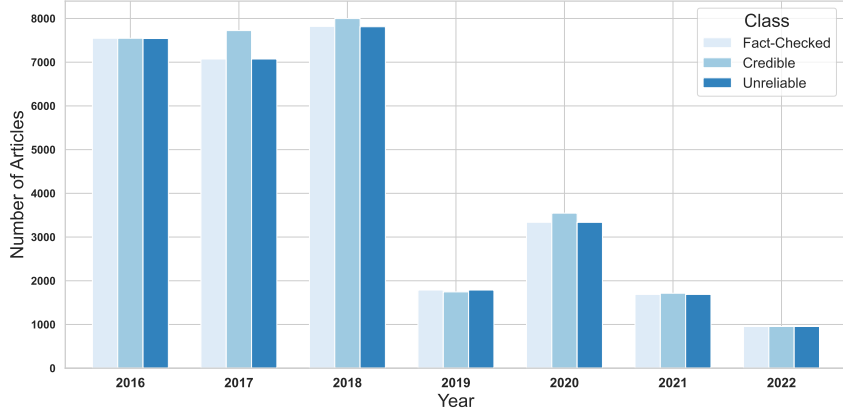


Figure 2: Number of articles per year in CREDULE.

Subsequently, the veracity prediction mechanism  $\mathcal{F}(\cdot)$  can be applied to the claim  $C$  with the filtered text evidence  $E'_T$  and the image evidence  $E_I$ :

$$y = \mathcal{F}(C, E'_T, E_I) \quad (2)$$

where  $y$  represents the predicted veracity label.  $\mathcal{G}(\cdot)$  serves as a critical component in mitigating the impact of trustworthiness issues and information leakage in AFC systems.

### 3.2 Constructing the CREDULE dataset

CREDULE comprises three distinct classes: Fact-Checked, Credible, and Unreliable, tailored for the classification of news articles to facilitate evidence filtering during the fact-checking process. Ensuring a balanced representation, each class encompasses a similar number of articles, consistent distributions in publication year, title length, and thematic content. Maintaining an even distribution of articles across years ensures comprehensive coverage of events over time. Entries in our dataset are presented in both article title and full-text formats whenever available.

To construct CREDULE, we merge information from six distinct sources: MultiFC [8], Politifact Fact-Check [9], PUBHEALTH [10], NELA-GT [11, 12, 13, 14, 15, 16], Fake News Corpus [17], and Getting Real About Fake News [18] datasets. These diverse datasets provide a rich and comprehensive foundation for our dataset creation process.

#### 3.2.1 Fact-checked class

In order to construct the fact-checked class, we leverage the following dataset: 1) **MultiFC**: encompassing 36,534 claims sourced from 24 Fact-checked domains, spanning until 2019. Each claim is accompanied by its respective veracity label, publication date, article title, URL, and additional metadata. 2) **PUBHEALTH**: comprising 11,832 claims within the health domain, sourced from various fact-checking websites, this dataset offers valuable insights. Alongside each claim, the dataset provides the article’s URL, publication date, and the main text body. 3) **Politifact**: We incorporate 21,152 statements from Politifact, spanning the years 2008 to 2022, sourced from the Politifact Fact Check dataset on Kaggle. Each entry includes URLs, truthfulness labels, and additional metadata. To ensure data consistency, we extract article titles from the URLs, omitting the term “Politifact” to mitigate biases.

From our collection of fact-checking articles, we retain only those published between 2016 and 2022. We also filter out duplicate and empty entries, resulting in 30,209 articles. Additionally, we employ a pre-trained DistilBert model, trained on the News Category Dataset<sup>6</sup>, to extract article topics based on their titles. This process yields a comprehensive set of 42 categories, which we further consolidate into 12 distinct groups for clarity. Furthermore, we extract the full text-body of articles where unavailable. First, we utilize the Beautiful Soup package to obtain the full article texts from the Politifact Fact Check dataset. Similarly, we develop domain-specific scripts to extract articles from the MultiFC dataset, excluding those already present in the Politifact dataset. The PUBHEALTH dataset provides full-text articles, obviating the need for extraction.

<sup>6</sup>HuggingFace Model: <https://huggingface.co/Yueh-Huan/news-category-classification-distilbert>

Table 1: CREDULE dataset statistics, constructed by combining and filtering various datasets.

Dataset	Class	Used Entries	Full-Texts
MultiFC	Fact-checked	16,057	70.5%
PUBHEALTH	Fact-checked	3,683	All
Politifact	Fact-checked	10,469	99.7%
Fake News Corpus	Credible	22,245	91.4%
NELA-GT	Credible+Unreliable	31,636	All
GRAFN	Unreliable	7,542	All

### 3.2.2 Credible and Unreliable Classes

To construct the Credible and Unreliable classes, we leverage the following datasets: 1) **Fake News Corpus:** It comprises millions of news articles, including both credible and non-credible sources. We specifically utilize the credible subset, extracted from reputable news outlets. To ensure data integrity, we extract article titles and dates directly from URLs, addressing issues with inaccuracies and duplicates. Additionally, we remove domain names from titles and employ the topic-extraction model to categorize the articles. 2) **NELA-GT Datasets:** Spanning from 2017 to 2022, these contain articles sourced from various domains, classified based on their credibility and bias. We extract pertinent information such as article titles, dates, domains, and topics using the topic-extraction model. These datasets are instrumental in constructing both the credible and non-credible classes. 3) **Getting Real about Fake News Dataset:** The Getting Real about Fake News dataset (GRAFN), sourced from Kaggle, features non-credible articles from 2016, scraped from 244 websites labeled as “bullshit” by the BS Detector Chrome Extension. We filter the dataset to include only English articles categorized as ‘bs’ (bullshit), ‘conspiracy’, ‘satire’, ‘junksci’, and ‘fake’. Obtaining article metadata, including titles, dates, domains, and main text bodies, enables us to further categorize articles based on their topics.

To ensure class balance across different years (2016-2022) and topics, we supplement articles from the Fake News Corpus and, if necessary, the NELA-GT datasets for the Credible class. We ensure an equal distribution of articles across various topics to align with the Fact-checked class, extracting the same number of articles per topic. Furthermore, for the Unreliable class, we incorporate articles from both NELA-GT and Getting Real About Fake News datasets. The distribution of articles across different years and classes, can be seen in Figure 2. CREDULE is balanced across classes in terms of articles per year and in terms of the length of titles within each class. We make use of the full text articles provided by NELA-GT and Getting Real About Fake News datasets. For the Fake News Corpus, we utilize the Beautiful Soup library and develop custom scripts to extract the content of full articles.

### 3.2.3 Domain Credibility Scores (DCS)

We augment CREDULE by integrating external domain credibility scores (DCS) from the Media Bias/Fact Check (MBFC) website. MBFC is an independent platform dedicated to combating media bias and misinformation and employs a rigorous evaluation combining objective metrics and subjective analysis to rate media sources based on factors such as bias, factual accuracy, and overall credibility. Bias assessments range from least biased to extreme bias on a scale of 0 to 10, while factuality scores vary from very high to very low based on fact-checking frequency and inclusion of critical information. Categorized into three tiers—high, medium, and low credibility—the final MBFC rating identifies highly credible sources with a score of 6 or above, medium credibility for scores ranging from 3 to 5, and low credibility for scores of 0 to 2 or sources rated as questionable, conspiracy, or pseudoscience. For each domain in CREDULE, we obtain DCS including ‘Bias Rating’, ‘Factual Reporting’ and ‘MBFC Credibility Rating’. This process involves querying the MBFC website for each domain, accessing the corresponding page, and extracting the relevant scores. We are able to extract domain credibility scores for 92.2% of the articles in CREDULE.

### 3.2.4 CREDULE Final Statistics

The CREDULE comprises a total of 91,632 articles, spanning the years 2016 to 2022 and classified into three distinct categories: Fact-checked (30,209), Credible (31,230), and Unreliable (30,193). Each article entry includes its title, publication date, URL, assigned topic, and classification label. Moreover, 92.7% of the articles in the dataset are accompanied by their full-text content. Notably, the dataset exhibits balanced distribution across all classes in terms of articles per year, topics, and title lengths. The statistics of the dataset are summarized in Table 1.

### 3.3 Evidence Verification Network

After constructing CREDULE, we develop an Evidence VERification Network (EVVER-Net) which can be used to discern between credible, unreliable and leaked evidence. EVVER-Net is a neural network classifier, which can be expressed as follows:

$$\hat{y}_i = \text{Softmax}(\mathbf{W}_1 \cdot \text{GELU}(\mathbf{W}_0 \cdot [T(e_i)[< EOS >]; s_i])) \quad (3)$$

where  $T(\cdot)$  stands for a Transformer backbone encoder,  $[< EOS >]$  for the position of the end-of-sentence (EOS) or classification token (CLS), depending on the encoder, of  $T(\cdot)$ ,  $s_i \in \mathbb{R}^1$  is the “domain credibility score” of  $e_i$ ,  $[\cdot]$  stands for concatenation,  $\mathbf{W}_0 \in \mathbb{R}^{1 \times h+1}$  is a GELU activated fully connected layer with  $h$  hidden dimensions and  $\mathbf{W}_1 \in \mathbb{R}^{h+1 \times 3}$  is the final classification layer activated with Softmax for 3 classes. Equation 3 represents the case where the network only has a single hidden layer  $l$ .

We develop three different versions of EVVER-Net. The first handles short texts (title articles), the second long texts (full articles), and the third also leverages domain credibility scores.

For “short text” experiments, we explore pre-trained Transformer-based language models such as T5 [21], DeBERTa [19] and CLIP [20]. For our experiments on full article texts, we employ transformer-based models tailored for processing larger text sequences. More specifically, we utilize Longformer [23], a model that integrates both local (window-based) and global attention mechanisms. For feature extraction, we utilize the [CLS] token to capture contextual information and train our classifier. Finally, we similarly employ LongT5 [22], an extension of the T5 model, suitable for long texts.

In order to incorporate DCS into EVVER-Net, we encode Factuality Scores categories ‘satire’(-3), ‘very low’(-2), ‘low’(-1), ‘mostly factual’(2), ‘high’(3) and ‘very high’(4). Articles without factuality scores are assigned a value of 0. For articles with a ‘mixed’ score, we consider the MBFC Credibility Rating. If it indicates ‘medium credibility’, ‘mixed’ is mapped to 1. Conversely, if it is ‘high credibility’ or ‘low credibility’, ‘mixed’ is mapped to 2 or -1, respectively. This differentiation ensures that we account for each case of ‘mixed’ credibility and the absence of data. Finally, we normalize DCS into a range of (0,1), concatenate them with the text embeddings and pass the combined input through EVVER-Net.

### 3.4 Implementation Details

We implement EVVER-Net using the PyTorch deep learning framework, leveraging its efficiency for training neural networks. To ensure reproducibility, we set the random seed to 42 before conducting any experiments. The dataset is divided into three subsets using an 80/10/10 training/validation/test split.

For both short- and long-text models, we utilize 3-fold Cross Validation and Grid Search, respectively to optimize and fine-tune EVVER-Net. In both cases, we explore various configurations including hidden sizes  $h \in \{512, 1024\}$  and number of layers  $l \in \{1, 2, 3\}$ . We employ the Adam optimizer with learning rates  $lr \in \{1e-3, 5e-4, 1e-4, 5e-5\}$  and batch sizes  $b \in \{512, 1024, 2048\}$ . Additionally, we experiment with dropout rates  $d \in \{0.1, 0.2, 0.25\}$  and L2 regularization  $r \in \{0, 1e-2, 1e-3\}$  to prevent overfitting and improve generalization.

Finally, we experiment with baseline models, namely Logistic Regression, Naive Bayes, Decision Trees and Multi-layer Perceptron (MLP) trained only on domain credibility scores or on statistical feature extraction methods like CountVectorizer and TF-IDF.

## 4 Results

### 4.1 Quantitative Results

Table 2 demonstrates the results of the baseline classifiers. We observe that Count Vectorizer achieves 66.3% accuracy with Logistic Regression and 69.4% with the MLP Classifier while TF-IDF yields 67.0% accuracy with Naive Bayes and 69.5% with an MLP classifier. When only leveraging DCS, a Decision Tree classifier reaches 68.5% accuracy, closely competing with text-based models.

Table 3 illustrates the performance of EVVER-Net while leveraging different transformer-based encoders. We observe that with short texts (titles) and without DCS, EVVER-Net reaches the best performance of 79.5% accuracy while employing embeddings from DeBERTa<sup>7</sup> and having hidden dimensions  $h = [512, 1024, 1024]$ , learning rate  $lr = 5e-5$ , dropout  $d = 0.1$ , l2 regulation of  $r = 1e-3$  and batch size  $b = 1024$ . This performance is followed by

<sup>7</sup><https://huggingface.co/microsoft/deberta-base>

Table 2: Baseline Experiments on CREDULE.

Input	Classifier	Accuracy
Count Vectorizer	Logistic Regression	66.3%
Count Vectorizer	MLP Classifier	69.4%
TF-IDF	Naive Bayes	67.0%
TF-IDF	MLP Classifier	69.5%
Domain Scores Only	Decision Trees	68.5%

Table 3: Performance of EVVER-Net on CREDULE for short or long texts, with different backbone encoder and with or without Domain Credibility Scores (DCS).

Encoder	Accuracy w/o DCS	Accuracy w/ DCS	Input
T5	78.3%	88.6%	short texts
Clip Text	79.3%	91.0%	short texts
DeBERTa	79.5%	91.5%	short texts
Long T5	79.1%	89.5%	long texts
Longformer	89.0%	94.4%	long texts

integrating CLIP’s text encoder <sup>8</sup> reaching 79.3% accuracy and then T5 <sup>9</sup> reaching 78.3% accuracy. When employing long texts (full articles), EVVER-Net with Longformer <sup>10</sup> as the backbone encoder, achieves 89% accuracy with  $h = [1024, 1024, 1024]$ ,  $lr = 5e - 4$ ,  $d = 0.2$ ,  $r = 1e - 3$  and  $b = 2048$ .

Furthermore, incorporating DCS from MBFC, can significantly and consistently enhance the classification accuracy of EVVER-Net across all 5 Transformer encoders. For short texts, the accuracy of EVVER-Net with DeBERTa increases to 91.5%, reflecting a notable relative improvement of +12% while for long texts, Longformer facilitates a substantial boost, reaching 94.4%, a +5.4% relative improvement. Overall, these results demonstrate notable improvements in classification accuracy across all backbone model encoder, highlighting the effectiveness of incorporating domain-specific characteristics into EVVER-Net.

## 4.2 Qualitative Analysis and Inference

In this section, we apply EVVER-Net to existing datasets, encompassing various text, multimodal, and Web-sourced datasets. We do not use DCS in this section because the domain names from which the evidence was collected are not provided by these datasets. By examining the classifier’s performance across diverse datasets, we aim to assess its robustness and applicability in different contexts. Specifically, we focus on datasets containing fact-checked articles, as well as those incorporating external evidence from various sources on the Web. Our evaluation begins with testing on text datasets, where the classifier’s ability to discern credible information from unreliable is put to the test. We then extend our analysis to multimodal datasets and finally, we explore datasets that aggregate evidence from the Web, mirroring the real-world application scenario of EVVER-Net. Results are summarized in Table 4

<sup>8</sup><https://huggingface.co/openai/clip-vit-base-patch32>

<sup>9</sup><https://huggingface.co/google-t5/t5-large>

<sup>10</sup><https://huggingface.co/allenai/longformer-base-4096>

Table 4: Inference Results: Applying EVVER-Net on evidence from various datasets.

Dataset	Data Type	Fact-checked	Credible	Unreliable	Samples
LIAR-PLUS	Ruling statements of fact-checked articles	98.5%	1.0%	0.5%	10,238
MOCHEG	Evidence snippets from fact-checked articles	83.6%	2.4%	14.0%	27,528
FACTIFY	Article from “Support” & “Insufficient” classes	6.4%	88.0%	5.6%	34,000
FACTIFY	Articles from “Refute” class	95.0%	3.5%	1.5%	8,500
NewsCLIPPings+	Article titles (Web)	14.1%	64.5%	21.4%	45,907
VERITE	Article titles (Web) - across 3 classes	35.3%	42.6%	22.1%	1,611
VERITE	Article titles (Web) - ‘True’ and ‘Miscaptioned’ classes	45.2%	35.2%	19.6%	1,094
VERITE	Article titles (Web) - ‘Out-of-Context’ class	14.5%	58.2%	27.3%	517



DATASET	CLAIM	EVIDENCE	EVVER-Net
LIAR-PLUS	"Just about everyone everywhere is spending more hours on the job, less time with their families, bringing home smaller and smaller paychecks, while they're paying more and more at the gas pump and the grocery stores" <b>Verdict: "Mostly True"</b>	"We did not assign credit to Kasich for his statement in March that Ohioans' wages have risen by more than \$10 billion since 2010. But we rated the statement as True. [...]" (source: Politifact)	Fact-Checked
MOCHEG	"President Eisenhower said that a political party must be dedicated to the advancement of a moral cause, otherwise it is just a conspiracy to seize power." <b>Verdict: "Correct Attribution"</b>	"If a political party does not have its foundation in the determination to advance a cause that is right and that is moral, then it is not a political party; it is merely a conspiracy to seize power." (source: Snopes)	Fact-Checked
NewsCLIPPings+	"Freeza Meats of Newry said it had been asked by an Irish company to store the meat after they had declined to buy it" <b>Verdict: "Pristine"</b>	"Food company Freeza Meats fined £70,000 costs" (source: BBC News)	Credible
VERITE	"An image shows a large crowd gathered on the Rio de Janeiro seafont to celebrate a Mass delivered by Pope Francis in 2013." <b>Verdict: "Truthful"</b>	"Pope Francis wraps up Brazil trip with Mass for 3 million" (source: CBS News)	Credible
	"A photograph captured by NASA's Mars Curiosity Rover on May 7, 2022, showed an artificial portal nearby." <b>Verdict: "Out-Of-Context"</b>	"Does This NASA Photo Show a 'Portal' and 'Wall' on Mars?" (source: Snopes)	Fact-Checked
	"A photograph shows plaster cast of a victim in Pompeii during the volcanic eruption in 79 A.D." <b>Verdict: "Truthful"</b>	"Time Traveling to Pompeii : r/TikTokCringe" (source: Reddit )	Unreliable
FACTIFY - "Support"	"Now that Sen. Kamala Harris is in, here are the Democrats who've said they're running for president." <b>Verdict: Support-Text</b>	"California Sen. Kamala Harris announced Monday that she will run for president in 2020, joining an increasingly crowded field of Democrats seeking to challenge President Donald Trump. [...] Booker and Sanders will also speak at a rally at South Carolina's state house, and Castro, the former mayor of San Antonio, is marching in the city's Martin Luther King Jr. Day parade." (source: ABC News)	Credible
FACTIFY - "Refute"	"Microsoft bought Sony for \$121 billion." <b>Verdict: Refute</b>	"A satire article claiming that Microsoft has bought Sony for \$130 billion has gone viral with some media outlets reporting it as true. BOOM found the original article in a Spanish website was intended to be a prank. [...] Furthermore, considering Microsoft and Sony are two of the biggest brands in the world, any news of one acquiring the other would make headlines. However, we found no news reports on Microsoft buying Sony." (source: BOOM Live)	Fact-Checked

Figure 3: Inference examples from EVVER-Net applied on the evidence of various datasets. We manually included the domain names from which the evidence originates, as they are not provided by the datasets.

#### 4.2.1 LIAR-PLUS Dataset

We initiate our evaluation with the LIAR-PLUS dataset [24], a well-known repository of fact-checked articles. LIAR-PLUS provides not only labeled articles but also detailed justifications for the veracity of each claim, which are sentences extracted from the section ‘Our Ruling’ of each fact-checked article. Therefore, we would expect high rates of “fact-checked” class. Indeed, when we extract features with DeBERTa and apply EVVER-Net on the “justification evidence”, it successfully recognizes fact-checked evidence pieces, with 98.5% being classified as Fact-checked, 1% as credible, and 0.5% as Unreliable. These findings underscore EVVER-Net ability to discern information from fact-check websites, affirming its utility in real-world fact-checking applications.

#### 4.2.2 MOCHEG Dataset

We further extend our evaluation by applying EVVER-Net to the MOCHEG dataset [26], which aggregates data from Politifact and Snopes websites. Specifically, we focus on the ‘evidence’ column of Corpus 2, comprising highlighted text snippets from fact-checked articles. Again, we observe consistent performance, with rate of 83.6% of evidence pieces being classified as fact-checks.

#### 4.2.3 FACTIFY Dataset

The next dataset we evaluate is the FACTIFY dataset [25]. This comprises textual and image claims, each associated with a reliable source of information referred to as a ‘document’. Claims are categorized into three classes: support, insufficient, and refute, based on their relationship with the corresponding document. Authors collected and scraped news articles from various credible sources in the US and India to compile data for the support and insufficient categories. We extract features with Longformer from the full-text article bodies and pass them through our classifier. Articles from the “Support” and “Insufficient” classes were collected from credible sources while articles in the “Refute” classes were collected from fact-checking websites. Our analysis reveal that 88% of the contents categorized as “Support” or “Insufficient” are classified as credible by the classifier, while only 6.4% and 5.6% are classified as Fact-checked and Unreliable, respectively. Notably, EVVER-Net correctly identifies that 95% of the documents in the “Refute” class as Fact-checked.

#### 4.2.4 NewsCLIPpings+ Dataset

We extend our evaluation to include the evidence sourced from [28], referred to as NewsCLIPpings+, which consists of external information collected from the Web. Abdelnabi et al. [28] utilized the NewsCLIPpings dataset [36], which contains both pristine and falsified (algorithmically decontextualized) image-caption pairs. In this setup, textual evidence was obtained by querying images in an inverse search mode using the Google Vision API. For our analysis, we utilize the test set of NewsCLIPpings+, comprising 7,264 image captions and 51,799 scraped text evidence pieces. After removing inaccurately scraped entries, such as image file names or ‘Page Not Found’ entries, we are left with 45,907 textual evidence pieces. EVVER-Net estimates that 64.5% of the evidence is Credible, 21.4% Unreliable, and 14.1% Fact-checked. Despite the majority of the evidence being considered credible, the presence of some pieces sourced from unreliable or leaked origins in NewsCLIPpings+ may not only compromise the accuracy of the veracity predictions but also introduce unrealistic elements into the assessment process.

#### 4.2.5 VERITE Dataset

We also assess evidence extracted in [41] for the VERITE dataset [27], which integrates image-caption pairs similar to the previous dataset. In this dataset, ‘MisCaptioned’ pairs sourced from fact-checked articles like Snopes and Reuters were considered. The authors collected misleading claims along with their associated images, utilizing the Google API to retrieve textual and visual evidence. After cleaning the gathered evidence, we obtain 1,611 text snippets, gathered from querying 1,000 image captions, classified as ‘True’, ‘Miscaptioned’, or ‘Out-of-Context’. Our classifier categorizes these snippets with 42.6% labeled as Credible, 35.3% as Fact-checked, and 22.1% as Unreliable. It is important to note that the higher percentage of Fact-checked snippets in this dataset is due to sourcing claims and images directly from fact-checked articles in the ‘Miscaptioned’ class. When these claims are searched using the Google API, it often returns the original or similar fact-checking articles, resulting in leaked evidence. For image-caption pairs classified as ‘Out-of-Context’, we observe a lower Fact-checked (14.5%) and a higher Unreliable percentage (27.3%) compared to the other categories. In this case, the claims are taken from fact-checked articles but the out-of-context images are retrieved from the Web. Thereafter, these images are used to retrieve textual information from a Google API thus increasing the likelihood of unreliable evidence.

#### 4.2.6 Inference

Our findings reveal that widely-used datasets incorporate leaked and unreliable evidence during the AFC process. This highlights the critical need for robust filtering mechanisms, like EVVER-Net, to identify and exclude such information effectively. Figure 3 provides examples sourced from LIAR-PLUS, MOCHEG, NewsCLIPpings+, VERITE, and FACTIFY datasets, alongside their classification by EVVER-Net. For instance, consider how “[...] But we rated the statement as True.” (LIAR-PLUS) or “A satire article [...] intended to be a prank [...]” (FACTIFY) - taken directly from within fact-checked articles - provide information curated by fact-checkers that directly support or refute the claim, thus creating an unrealistic scenario for early detection of new misinformation. On the other hand, short article titles such as “Does This NASA Photo Show a ‘Portal’ and ‘Wall’ on Mars?” or “Time Traveling to Pompeii: r/TikTokCringe” (VERITE), does not necessarily provide any unreliable or leaked information by themselves. However, if the articles’ content were to be collected and analysed, they would be problematic for AFC systems.

## 5 Conclusion

In this study, we address the critical but overlooked challenges associated with verifying the quality of external information used as evidence in existing datasets, particularly the presence of leaked and unreliable information. We develop the CREDULE dataset, comprising 91,632 news articles classified as Fact-checked, Credible, and Unreliable. Additionally, we introduce the Evidence Verification Network (EVVER-Net), a robust solution for evidence verification and filtering, and conducted experiments with various language models. EVVER-Net reaches 79.5% and 89.0% accuracy for short and long texts, respectively, without utilizing domain credibility scores. By utilizing domain credibility scores, the performance of EVVER-Net further improves to 91.5% and 94.4%. Furthermore, our analysis on widely used datasets, including LIAR-PLUS, MOCHEG, FACTIFY, NewsCLIPpings+ and VERITE, reveals that collected evidence often contain leaked and unreliable information, thereby diminishing the effectiveness and realism of building robust AFC systems. These findings underscore the importance of implementing evidence verification and filtering solutions such as EVVER-Net during the evidence retrieval task of AFC systems.

While our study offers valuable insights into evidence verification in AFC, it has certain limitations. First, it solely focuses on textual evidence, overlooking the potential influence of images and videos. Visual elements can also contain unreliable information (e.g., by being fabricated, manipulated or synthetically generated) or having artifacts that indirectly indicate that they are sourced from fact-checked articles (e.g., watermarks). Future research should aim to

enhance the robustness of evidence verification and filtering methods like EVVER-Net by expanding the CREDULE dataset to include a more diverse range of articles and sources, including multimedia content. Furthermore, we refrain from collecting new external information from the Web and filtering them with EVVER-Net in order to examine its impact on AFC systems. We hypothesize that the performance of AFC systems would decrease if they previously primarily relied on leaked evidence, which would afterwards be removed by EVVER-Net. On the other hand, we hypothesize that removing unreliable information would improve performance, thus somewhat balancing out the decrease in performance from removing leaked evidence. Overall, this would result into a more realistic and reliable framework for training and evaluating AFC systems, especially on new and emerging misinformation. Nevertheless, future research should systematically examine these effects.

## Acknowledgements

This work is partially funded by the Horizon Europe projects vera.ai under grant agreement no. 101070093 and DisAI under grant agreement no. 101079164.

## References

- [1] Andrew Duffy, Edson Tandoc, and Rich Ling. Too good to be true, too good not to share: the social utility of fake news. *Information, Communication & Society*, 23(13):1965–1979, 2020.
- [2] Femi Olan, Uchitha Jayawickrama, Emmanuel Ogiemwonyi Arakpogun, Jana Suklan, and Shaofeng Liu. Fake news on social media: the impact on society. *Information Systems Frontiers*, pages 1–16, 2022.
- [3] Ray Oshikawa, Jing Qian, and William Yang Wang. A survey on natural language processing for fake news detection. *arXiv preprint arXiv:1811.00770*, 2018.
- [4] Hanqing Zhao, Wenbo Zhou, Dongdong Chen, Tianyi Wei, Weiming Zhang, and Nenghai Yu. Multi-attentional deepfake detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2185–2194, 2021.
- [5] Firoj Alam, Stefano Cresci, Tanmoy Chakraborty, Fabrizio Silvestri, Dimiter Dimitrov, Giovanni Da San Martino, Shaden Shaar, Hamed Firooz, and Preslav Nakov. A survey on multimodal disinformation detection. *arXiv preprint arXiv:2103.12541*, 2021.
- [6] Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. A survey on automated fact-checking. *Transactions of the Association for Computational Linguistics*, 10:178–206, 2022.
- [7] Max Glockner, Yufang Hou, and Iryna Gurevych. Missing counter-evidence renders NLP fact-checking unrealistic for misinformation. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5916–5936, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.
- [8] Isabelle Augenstein, Christina Lioma, Dongsheng Wang, Lucas Chaves Lima, Casper Hansen, Christian Hansen, and Jakob Grue Simonsen. Multifc: A real-world multi-domain dataset for evidence-based fact checking of claims. *arXiv preprint arXiv:1909.03242*, 2019.
- [9] Rishabh Misra. Politifact Fact Check Dataset. <https://www.kaggle.com/datasets/rmisra/politifact-fact-check-dataset>, 2022.
- [10] Neema Kotonya and Francesca Toni. Explainable automated fact-checking for public health claims. *arXiv preprint arXiv:2010.09926*, 2020.
- [11] Benjamin Horne, Sara Khedr, and Sibel Adali. Sampling the news producers: A large news and feature data set for the study of the complex media landscape. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 12, 2018.
- [12] Jeppe Nørregaard, Benjamin D Horne, and Sibel Adali. Nela-gt-2018: A large multi-labelled news dataset for the study of misinformation in news articles. In *Proceedings of the international AAAI conference on web and social media*, volume 13, pages 630–638, 2019.
- [13] Maurício Gruppi, Benjamin D. Horne, and Sibel Adali. Nela-gt-2019: A large multi-labelled news dataset for the study of misinformation in news articles, 2020.
- [14] Maurício Gruppi, Benjamin D Horne, and Sibel Adali. Nela-gt-2020: A large multi-labelled news dataset for the study of misinformation in news articles. *arXiv preprint arXiv:2102.04567*, 2021.

- [15] Maurício Gruppi, Benjamin D. Horne, and Sibel Adali. Nela-gt-2021: A large multi-labelled news dataset for the study of misinformation in news articles, 2021.
- [16] Maurício Gruppi, Benjamin D. Horne, and Sibel Adali. Nela-gt-2022: A large multi-labelled news dataset for the study of misinformation in news articles, 2023.
- [17] Maciej Szpakowski. FakeNewsCorpus Dataset. <https://github.com/several27/FakeNewsCorpus>, 2020.
- [18] Megan Risdal. Getting real about fake news, 2016.
- [19] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*, 2020.
- [20] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [21] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- [22] Mandy Guo, Joshua Ainslie, David Uthus, Santiago Ontanon, Jianmo Ni, Yun-Hsuan Sung, and Yinfei Yang. Longt5: Efficient text-to-text transformer for long sequences. *arXiv preprint arXiv:2112.07916*, 2021.
- [23] Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- [24] Tariq Alhindi, Savvas Petridis, and Smaranda Muresan. Where is your evidence: Improving fact-checking by justification modeling. In James Thorne, Andreas Vlachos, Oana Cocarascu, Christos Christodoulopoulos, and Arpit Mittal, editors, *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*, pages 85–90, Brussels, Belgium, November 2018. Association for Computational Linguistics.
- [25] Shreyash Mishra, S Suryavardan, Amrit Bhaskar, Parul Chopra, Aishwarya Reganti, Parth Patwa, Amitava Das, Tanmoy Chakraborty, Amit Sheth, Asif Ekbal, et al. Factify: A multi-modal fact verification dataset. In *Proceedings of the First Workshop on Multimodal Fact-Checking and Hate Speech Detection (DE-FACTIFY)*, 2022.
- [26] Barry Menglong Yao, Aditya Shah, Lichao Sun, Jin-Hee Cho, and Lifu Huang. End-to-end multimodal fact-checking and explanation generation: A challenging dataset and models. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2733–2743, 2023.
- [27] Stefanos-Iordanis Papadopoulos, Christos Koutlis, Symeon Papadopoulos, and Panagiotis C Petrantonakis. Verite: a robust benchmark for multimodal misinformation detection accounting for unimodal bias. *International Journal of Multimedia Information Retrieval*, 13(1):4, 2024.
- [28] Sahar Abdelnabi, Rakibul Hasan, and Mario Fritz. Open-domain, content-based, multi-modal fact-checking of out-of-context images via online resources. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14940–14949, 2022.
- [29] James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. FEVER: a large-scale dataset for fact extraction and VERification. In Marilyn Walker, Heng Ji, and Amanda Stent, editors, *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 809–819, New Orleans, Louisiana, June 2018. Association for Computational Linguistics.
- [30] Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dongwon Lee, and Huan Liu. Fakenewsnet: A data repository with news content, social context, and spatiotemporal information for studying fake news on social media. *Big data*, 8(3):171–188, 2020.
- [31] Kashif Khan, Ruizhe Wang, and Pascal Poupart. Watclaimcheck: A new dataset for claim entailment and inference. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1293–1304, 2022.
- [32] Naemul Hassan, Gensheng Zhang, Fatma Arslan, Josue Caraballo, Damian Jimenez, Siddhant Gawsane, Shohedul Hasan, Minumol Joseph, Aaditya Kulkarni, Anil Kumar Nayak, et al. Claimbuster: The first-ever end-to-end fact-checking system. *Proceedings of the VLDB Endowment*, 10(12):1945–1948, 2017.
- [33] Hannah Rashkin, Eunsol Choi, Jin Yea Jang, Svitlana Volkova, and Yejin Choi. Truth of varying shades: Analyzing language in fake news and political fact-checking. In *Proceedings of the 2017 conference on empirical methods in natural language processing*, pages 2931–2937, 2017.

- [34] Christina Boididou, Stuart E Middleton, Zhiwei Jin, Symeon Papadopoulos, Duc-Tien Dang-Nguyen, Giulia Boato, and Yiannis Kompatsiaris. Verifying information with multimedia content on twitter: a comparative study of automated approaches. *Multimedia tools and applications*, 77:15545–15571, 2018.
- [35] Shivangi Aneja, Chris Bregler, and Matthias Nießner. Cosmos: Catching out-of-context misinformation with self-supervised learning. *arXiv preprint arXiv:2101.06278*, 2021.
- [36] Grace Luo, Trevor Darrell, and Anna Rohrbach. Newsclippings: Automatic generation of out-of-context multimodal media. *arXiv preprint arXiv:2104.05893*, 2021.
- [37] Fuxiao Liu, Yinghan Wang, Tianlu Wang, and Vicente Ordonez. Visual news: Benchmark and challenges in news image captioning. *arXiv preprint arXiv:2010.03743*, 2020.
- [38] Kashyap Papat, Subhabrata Mukherjee, Jannik Strötgen, and Gerhard Weikum. Credibility assessment of textual claims on the web. In *Proceedings of the 25th ACM international on conference on information and knowledge management*, pages 2173–2178, 2016.
- [39] Kashyap Papat, Subhabrata Mukherjee, Andrew Yates, and Gerhard Weikum. Declare: Debunking fake news and false claims using evidence-aware deep learning. *arXiv preprint arXiv:1809.06416*, 2018.
- [40] Chris Samarinas, Wynne Hsu, and Mong Li Lee. Latent retrieval for large-scale fact-checking and question answering with nli training. In *2020 IEEE 32nd International Conference on Tools with Artificial Intelligence (ICTAI)*, pages 941–948, 2020.
- [41] Stefanos-Iordanis Papadopoulos, Christos Koutlis, Symeon Papadopoulos, and Panagiotis C Petrantonakis. Red-dot: Multimodal fact-checking via relevant evidence detection. *arXiv preprint arXiv:2311.09939*, 2023.
- [42] Georgi Karadzhov, Preslav Nakov, Lluís Màrquez, Alberto Barrón-Cedeño, and Ivan Koychev. Fully automated fact checking using external sources. *arXiv preprint arXiv:1710.00341*, 2017.
- [43] Michael Schlichtkrull, Zhijiang Guo, and Andreas Vlachos. Averitec: A dataset for real-world claim verification with evidence from the web. *Advances in Neural Information Processing Systems*, 36, 2024.